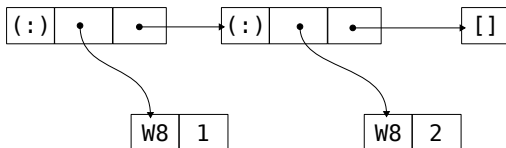


Hatékony Input/Output

A listák nem eléggé hatékonyak...

A szabványos listák nem ideálisak input/outputra:

- ▶ Sok a mutató, az ún. *cache miss*
- ▶ Sok helyet foglalnak, nagy az ún. *overhead*



Miként is fog ez pontosan működni Haskellben?

example = **do**

contents ← *readFile* "/etc/hosts"

return \$ *head* \$ *lines* *contents*

A readFile működése

```
readFile :: FilePath → IO String  
readFile name = openFile name ReadMode >>= hGetContents
```

```
hGetContents :: Handle → IO String  
hGetContents handle =  
  wantReadableHandle "hGetContents" handle $ \ h . do  
    xs ← lazyRead handle  
    return (h { haType = SemiClosedHandle }, xs)
```

-- unsafeInterleaveIO: az IO-akció csak az érték hivatkozásakor

-- (késleltetten) hajtódik végre

```
lazyRead :: Handle → IO String
```

```
lazyRead handle = unsafeInterleaveIO $
```

```
  withHandle "hGetContents" handle $ \ h . do
```

```
    case (haType h) of
```

```
      SemiClosedHandle → lazyReadBuffered handle h
```

```
      ClosedHandle     → ioException $
```

```
        IOError (Just handle) IllegalOperation "hGetContents"  
          "delayed read on closed handle" Nothing Nothing
```

```
    _ → ioException $
```

```
      IOError (Just handle) IllegalOperation "hGetContents"  
        "illegal handle type" Nothing Nothing
```

Data.ByteString

`http://hackage.haskell.org/package/bytestring`

Idő- és tárhatékony adatszerkezet, byte-ok véges (*finite sequence*) és végtelen sorozatához (*infinite stream*).

- ▶ *Mohó (strict) ByteString*. Egyetlen összefüggő memóriaterület, tömb. C és Haskell közti kommunikáció esetén kényelmes.
- ▶ *Lusta (lazy) ByteString*. Nagyobb összefüggő memóriaterületek lustán láncolt listája. Input/output adatfolyamok esetén hasznos.

import qualified ByteString as B

```
example = do
  contents ← B.readFile "/etc/hosts"
  return $ B.head $ B.lines contents
```

ByteString: Mohó változat (definíció)

```
{-# LANGUAGE BangPatterns #-}
```

```
import Foreign
```

```
-- minden ByteString költsége: 9 gépi szó
```

```
data ByteString = BS
```

```
    {-# UNPACK #-} !(ForeignPtr Word8) -- tartalom (puffer)
```

```
    {-# UNPACK #-} !Int                -- kezdőpozíció (offset)
```

```
    {-# UNPACK #-} !Int                -- hossz
```

```
empty :: ByteString
```

```
empty = BS nullForeignPtr 0 0
```

```
singleton :: Word8 → ByteString
```

```
singleton c = unsafeCreate 1 $ \ b . poke b c
```

```
unsafeCreate :: Int → (Ptr Word8 → IO ()) → ByteString
```

```
unsafeCreate l f = unsafePerformIO $ create l f
```

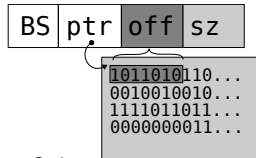
```
create :: Int → (Ptr Word8 → IO ()) → IO ByteString
```

```
create l f = do
```

```
    fp ← mallocForeignPtrBytes l
```

```
    withForeignPtr fp $ \ p . f p
```

```
    return $! BS fp 0 l
```



ByteString: Mohó változat (definíció)

```
pack :: [Word8] → ByteString
pack ws = unsafeCreate (length ws) $ \p . f p ws
  where
    f !_ []      = return ()
    f !p (x : xs) = poke p x >> f (p `plusPtr` 1) xs

unpack :: ByteString → [Word8]
unpack bs = unpackAppend bs []

unpackAppend :: ByteString → [Word8] → [Word8]
unpackAppend (BS fp off len) xs
  | len ≤ 100 = unpackAppendS (BS fp off len) xs
  | otherwise = unpackAppendS (BS fp off 100) remainder
  where
    remainder = unpackAppend (BS fp (off + 100) (len - 100)) xs

unpackAppendS :: ByteString → [Word8] → [Word8]
unpackAppendS (BS fp off len) xs = unsafePerformIO $ withForeignPtr fp $
  \base . loop (base `plusPtr` (off - 1))
    (base `plusPtr` (off - 1 + len)) xs
  where
    loop !sentinal !p acc
      | p ≡ sentinal = return acc
      | otherwise    = do
        x ← peek p
        loop sentinal (p `plusPtr` (-1)) (x : acc)
```

ByteString: Mohó változat (műveletek)

```
{-# LANGUAGE OverloadedStrings #-}
```

```
class IsString  $\alpha$  where
```

```
  fromString :: String  $\rightarrow$   $\alpha$ 
```

```
class Show  $\alpha$  where
```

```
  show ::  $\alpha$   $\rightarrow$  String
```

```
instance IsString ByteString where fromString = packChars
```

```
instance Show ByteString where show = unpackChars
```

```
take :: Int  $\rightarrow$  ByteString  $\rightarrow$  ByteString
```

```
take n bs@(BS ptr off len)
```

```
  | n  $\leq$  0    = empty
```

```
  | n  $\geq$  len = bs
```

```
  | otherwise = BS ptr off n
```

```
append :: ByteString  $\rightarrow$  ByteString  $\rightarrow$  ByteString
```

```
append (BS __ 0)      y                = y
```

```
append x              (BS __ 0)        = x
```

```
append (BS ptr1 off1 len1) (BS ptr2 off2 len2) =
```

```
  unsafeCreate (len1 + len2) $  $\lambda$  destptr1 . do
```

```
    let destptr2 = destptr1 `plusPtr` len1
```

```
    withForeignPtr ptr1 $  $\lambda$  p2 . memcpy destptr1 (p1 `plusPtr` off1) len1
```

```
    withForeignPtr ptr2 $  $\lambda$  p2 . memcpy destptr2 (p2 `plusPtr` off2) len2
```

ByteString: Lusta változat

Mohó (egyenként legfeljebb 64 KB méretű) *ByteString* értékek lustán láncolt listája.

```
data LazyByteString
    = Empty
  | Chunk !ByteString LazyByteString
```

- ▶ Byte-ok tiszta folyamának hatékony és tömör ábrázolási módja. A nem hivatkozott szeleteket a szeméthyűjtőgető felszabadítja.
- ▶ Az adatfolyamok feldolgozása konstans memóriahasználattal jár, de csak egyszerű, rövidéletű programok esetében javasolt.
- ▶ A szeletek határán kisebb a hatékonyság.
- ▶ A hibákat nem lehet egyértelműen jelezni.

ByteString: Lusta változat (műveletek)

```
import qualified ByteString as BS
```

```
empty :: LazyByteString  
empty = Empty
```

```
singleton :: Word8 → LazyByteString  
singleton w = Chunk (BS.singleton w) Empty
```

```
chunk :: ByteString → LazyByteString → LazyByteString  
chunk c@(BS.BS__len) cs  
  | len ≡ 0    = cs  
  | otherwise = Chunk c cs
```

```
pack :: [Word8] → LazyByteString  
pack ws = f 32 ws
```

where

```
  f n cs = case (BS.packUptoLenBytes n cs) of  
    (bs, []) → chunk bs Empty  
    (bs, cs') → Chunk bs $ f ((n * 2) `min` optimalChunkSize) cs'
```

```
  optimalChunkSize = 4 * 1024
```

```
unpack :: LazyByteString → [Word8]  
unpack Empty      = []  
unpack (Chunk c cs) = BS.unpackAppend c (unpack cs)
```

Data.ByteString.Builder

Az ismételt *append* nagyon költséges. Ezért a lusta *ByteString* értékekből készítsünk egy differencialistát:

- ▶ konstans idejű *append* műveletet kapunk
- ▶ rögzített vagy akár változó byte-hosszúságú értékeket tudunk *ByteString* típusúra alakítani
- ▶ alkalmazástól függő módon állíthatunk elő *ByteString* értéket

```
import qualified Data.ByteString.Lazy as L
import Data.ByteString.Builder
import Data.Foldable (foldMap)
import Data.Monoid
```

```
renderString :: String → Builder
renderString cs = charUtf8 "'" ⊕ foldMap escape cs ⊕ charUtf8 "'"
  where
    escape '\\\' = charUtf8 '\\\' ⊕ charUtf8 '\\\'
    escape '\"'  = charUtf8 '\\\' ⊕ charUtf8 '\"'
    escape c     = charUtf8 c

encodeUtf8 :: String → L.ByteString
encodeUtf8 = toLazyByteString o renderString
```

[10..18]

Data.ByteString.Builder: Implementáció (részlet)

newtype Builder = Builder ($\forall \rho . \text{BuildStep } \rho \rightarrow \text{BuildStep } \rho$)

-- BufferRange: Mutatók által megadott terület a memóriában.

type BuildStep α = BufferRange \rightarrow IO (BuildSignal α)

data BuildSignal α

= Done

{-# UNPACK #-} !(Ptr Word8) -- mutató az adat után
 α -- állapot

| BufferFull -- megtelt puffer jelzése
{-# UNPACK #-} !Int -- minimum szükséges méret
{-# UNPACK #-} !(Ptr Word8) -- mutató az adat után
(BuildStep α) -- következő lépés

| InsertChunk -- darab beillesztése
{-# UNPACK #-} !(Ptr Word8) -- mutató az adat után
{-# UNPACK #-} ByteString -- a beillesztendő darab
(BuildStep α) -- következő lépés

append :: Builder \rightarrow Builder \rightarrow Builder

append (Builder b_1) (Builder b_2) = Builder \$ $b_1 \circ b_2$

toLazyByteString :: Builder \rightarrow LazyByteString

toLazyByteString = toLazyByteStringWith

(safeStrategy smallChunkSize {- 4 KB -} defaultChunkSize {- 32 KB -}) Empty

Egy hatékonyabb megoldás: *Iteratee*

Az "iteratee" egy *kompozicionális* absztrakció, amellyel szekvenciálisan gyártott bemenetet tudunk lépésről lépésre feldolgozni.

Háromféle értékét kaphat:

- ▶ A bemenet következő szelete
- ▶ „Nincs több adat”
- ▶ „Iterálás vége”

Háromféle értéket (állapotot) adhat vissza (a hívó felé):

- ▶ „Megállás egy végeredménnyel”
- ▶ „Feldolgozás folytatása”
- ▶ „Hiba történt”

Iteratee: Definíció

```
import Control.Exception
import qualified Data.ByteString as BS
import qualified Data.ByteString.Lazy.Char8 as Lazy
```

```
data Chunk = Chunk
  { chunkData    :: !Lazy.ByteString
  , chunkAtEOF   :: !Bool
  }
```

```
newtype Iteratee  $\alpha$  = Iteratee (Chunk  $\rightarrow$  Result  $\alpha$ )
```

```
data Result  $\alpha$ 
  = Done { rResult ::  $\alpha$ , rResidual :: Chunk }
  | NeedInput !(Iteratee  $\alpha$ )
  | NeedIO !(IO (Result  $\alpha$ ))
  | Failed !SomeException
```

Iteratee: Sorok beolvasása (példa)

```
consumeLine :: Iteratee (Maybe Lazy.ByteString)
consumeLine = Iteratee (f Lazy.empty)
```

```
f :: Lazy.ByteString → Chunk → Result (Maybe Lazy.ByteString)
f acc (Chunk { chunkData = input, chunkAtEOF = eof })
```

```
  | ¬ (Lazy.null b) =
    Done { rResult    = Just acca
          , rResidual = Chunk { chunkData    = Lazy.tail b
                              , chunkAtEOF = eof
                              }
          }
    | ¬ eof = NeedInput (Iteratee (f acca))
    | otherwise =
```

```
    Done { rResult    = Nothing
          , rResidual = Chunk { chunkData    = acca
                              , chunkAtEOF = eof
                              }
          }
  }
```

where

```
(a, b) = Lazy.break (≡ '\n') input
acca  = acc `Lazy.append` a
```

Enumerator: Állomány feldolgozása (példa)

type *Enumerator* $\alpha = \text{Iteratee } \alpha \rightarrow \text{IO } (\text{Result } \alpha)$

enumerateFile :: *FilePath* \rightarrow *Enumerator* α

enumerateFile path = λ iter .

bracket (*openFile path ReadMode*) *hClose* (*g iter*)

g :: *Iteratee* $\alpha \rightarrow \text{Handle} \rightarrow \text{IO } (\text{Result } \alpha)$

g (*Iteratee iter*) *h* = **do**

input \leftarrow *BS.hGetSome h 32752*

if (*BS.null input*)

then *return* \$ *NeedInput* (*Iteratee iter*)

else *run* \$

iter (*Chunk* { *chunkData* = *Lazy.fromChunks* [*input*]
 , *chunkAtEOF* = *False*
 })

where

run (*NeedInput i*) = *g i h*

run (*NeedIO i*) = *i >>= run*

run result = *return result*

Iteratee: Minden együtt

```
chunkEOF :: Chunk  
chunkEOF = Chunk { chunkData = Lazy.empty, chunkAtEOF = True }
```

```
getResult :: Result  $\alpha$   $\rightarrow$  IO  $\alpha$   
getResult (Done { rResult = x }) = return x  
getResult (NeedInput (Iteratee f)) = getResult (f chunkEOF)  
getResult (NeedIO io)              = io >>= getResult  
getResult (Failed e)                = throwIO e
```

```
example = do  
  result  $\leftarrow$  enumerateFile "/etc/hosts" consumeLine  
  getResult result
```


Text: ByteString szövegekre

`http://hackage.haskell.org/package/text`

Unicode karaktersorozatok és azok kódolásainak hatékony, tömör, *immutable* (mohó és lusta) ábrázolása, komoly *loop fusion* támogatással.

```
import Data.Text as T
import qualified Data.Text.IO as T
import Data.Text.Encoding as E
import Data.ByteString (ByteString)

-- „subject to fusion”: köztes adatszerkezetek elhagyása (hatékony)
countChars :: ByteString → Int
countChars = T.length ∘ T.toUpper ∘ E.decodeUtf8

example = do
  contents ← T.readFile "/etc/hosts"
  return $ T.head $ T.lines contents
```

Text: Implementáció (bepillantás)

$$length :: Text \rightarrow Int$$
$$length = streamLength \circ stream$$
$$\mathbf{data} \text{ Step } \sigma \alpha = Done \mid Skip !\sigma \mid Yield !\alpha !\sigma$$

data Stream α = ∀ σ . Stream (σ → Step σ α) -- léptető függvény
!σ -- aktuális állapot

$$stream :: Text \rightarrow Stream Char$$

```
stream (Text arr off len) =
```

```
Stream next off (betweenSize (len `shiftR` 1) len)
```

where

$$!end = off + len$$

next !i

 $| i > \text{end} = \text{Done}$
$$0xD800 \leq n_1 \leq 0xDBFF = \text{Yield}(U16.chr2\ n_1\ n_2)(i + 2)$$
$$| otherwise = Yield (unsafeChr n_1) (i + 1)$$

where

$$[n_1, n_2] = \text{map } (A.\text{unsafeIndex arr}) [i, i + 1]$$
$$streamLength :: Integral \alpha \Rightarrow Stream Char \rightarrow \alpha$$
$$\text{streamLength}(\text{Stream next } s_0 \text{ len}) = \text{loop length } 0 \ s_0$$

where

```
loop_length!z s = case (next s) of
```

Done $\rightarrow z$

$$\text{Skip } s' \rightarrow \text{loop length } z \ s'$$
$$\text{Yield } s' \rightarrow \text{loop_length} (z + 1) s'$$

[18..18]